



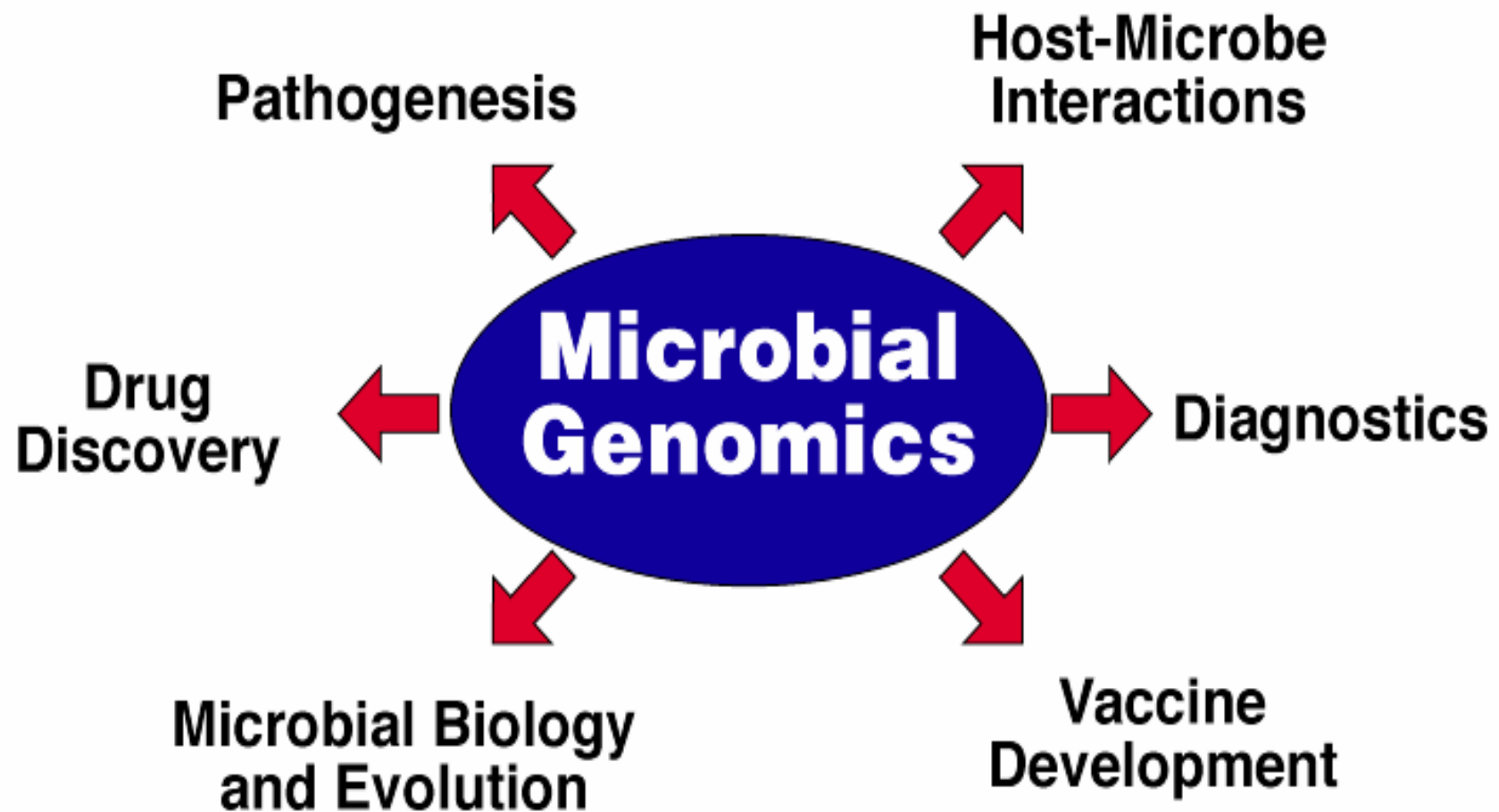
TIGR's MSC Program

Dr. Eric Eisenstadt, Ph.D.
VP for Research

Topics

(not necessarily in this order)

- MSC Overview and Mandate
- Genomes and Workflow across Resource Centers
- Data Release Policy & Reporting
- New Project Evaluations
- Current MSC Status & Completion
- Program Interactions



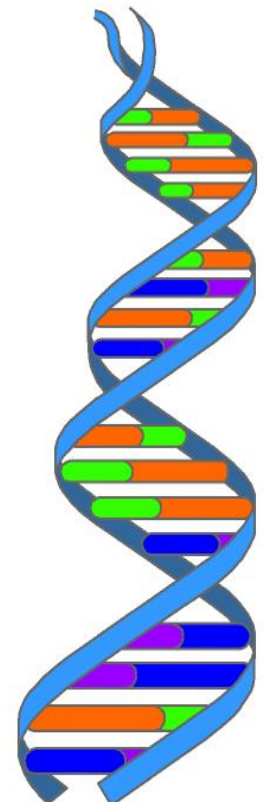
NIAID Microbial Genome Sequencing Centers

Goal:

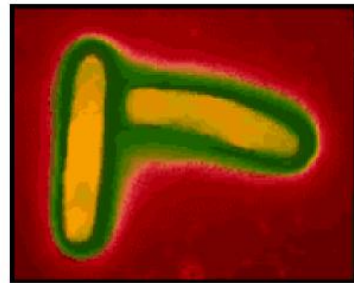
Rapid and cost-efficient production of high-quality genome sequences of human pathogens and invertebrate vectors of disease.

Features:

- Two Genome Sequencing Centers
 - The Institute for Genomic Research (TIGR)
 - MIT
- Capacity to sequence genomes for:
 - Other government agencies
 - Scientific community
 - Response to national emergencies



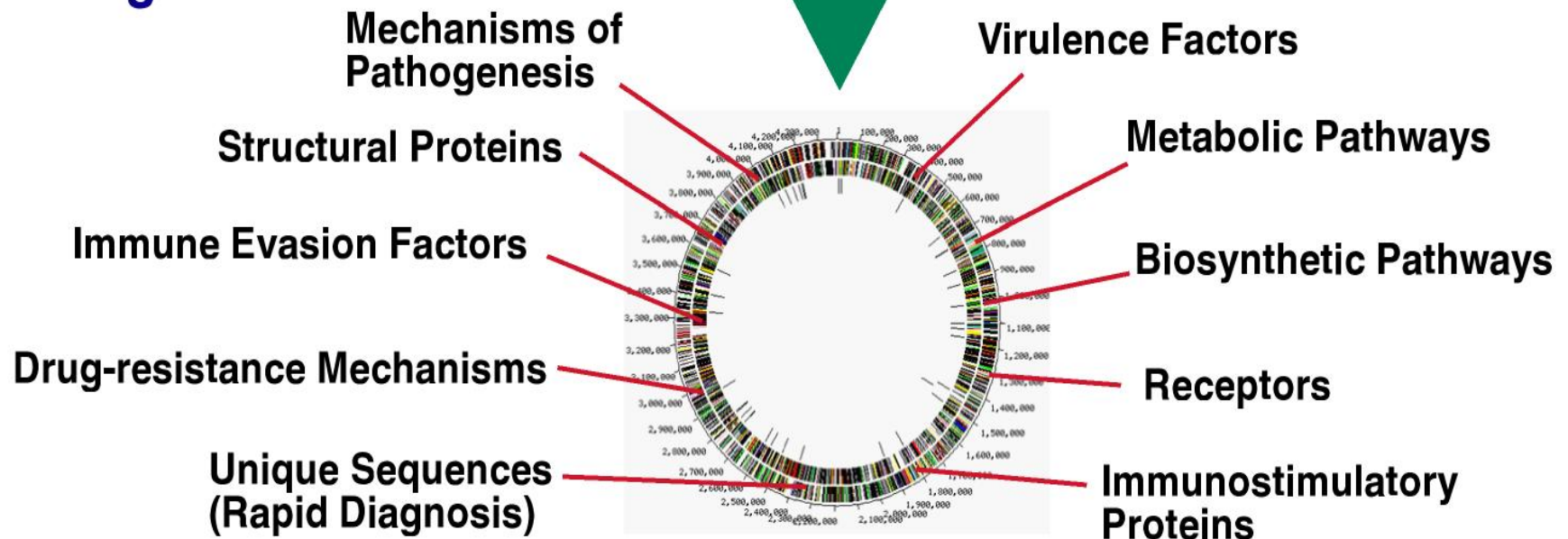
Mining Microbial Genomes to Develop Countermeasures



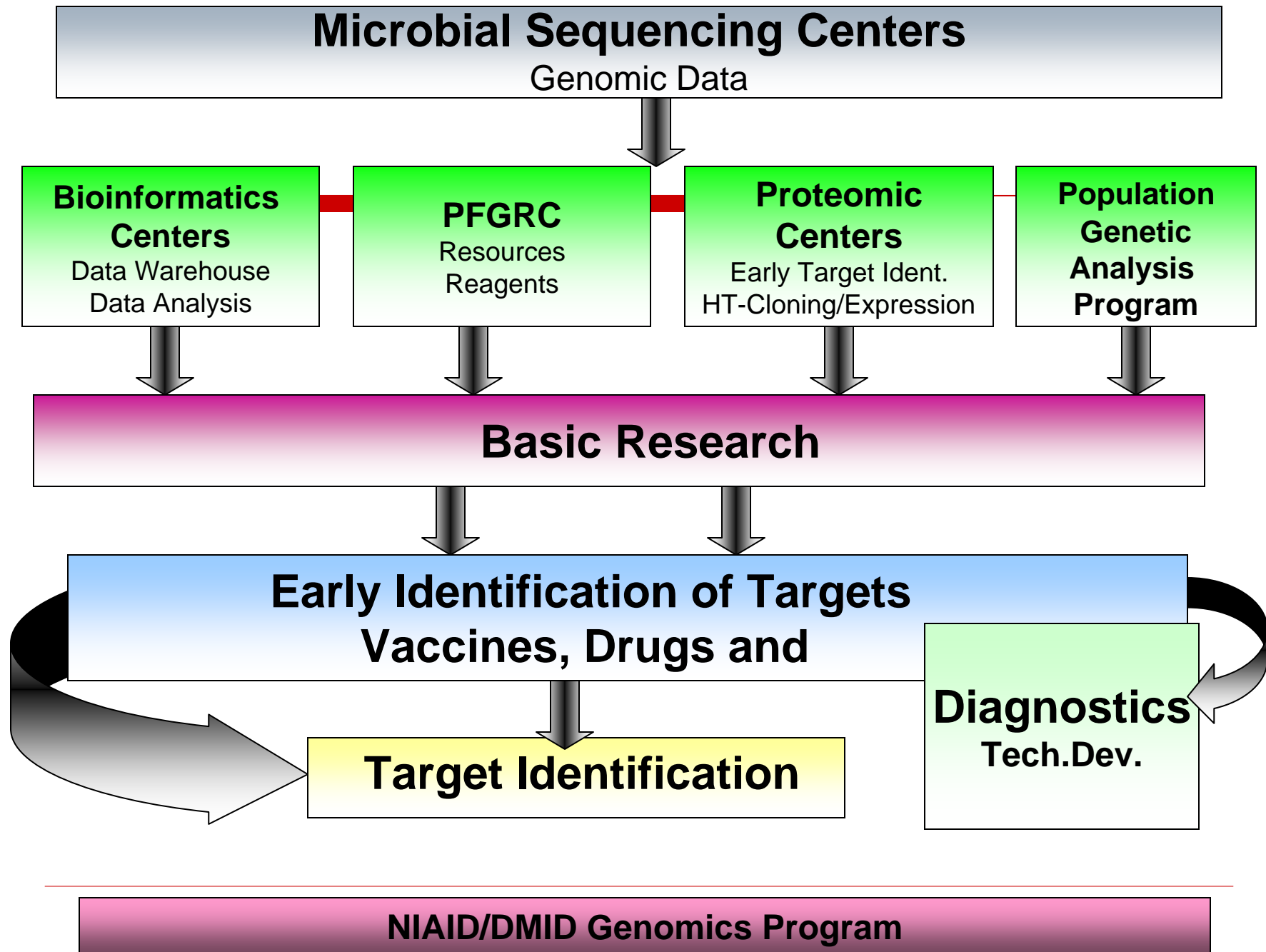
Pathogen



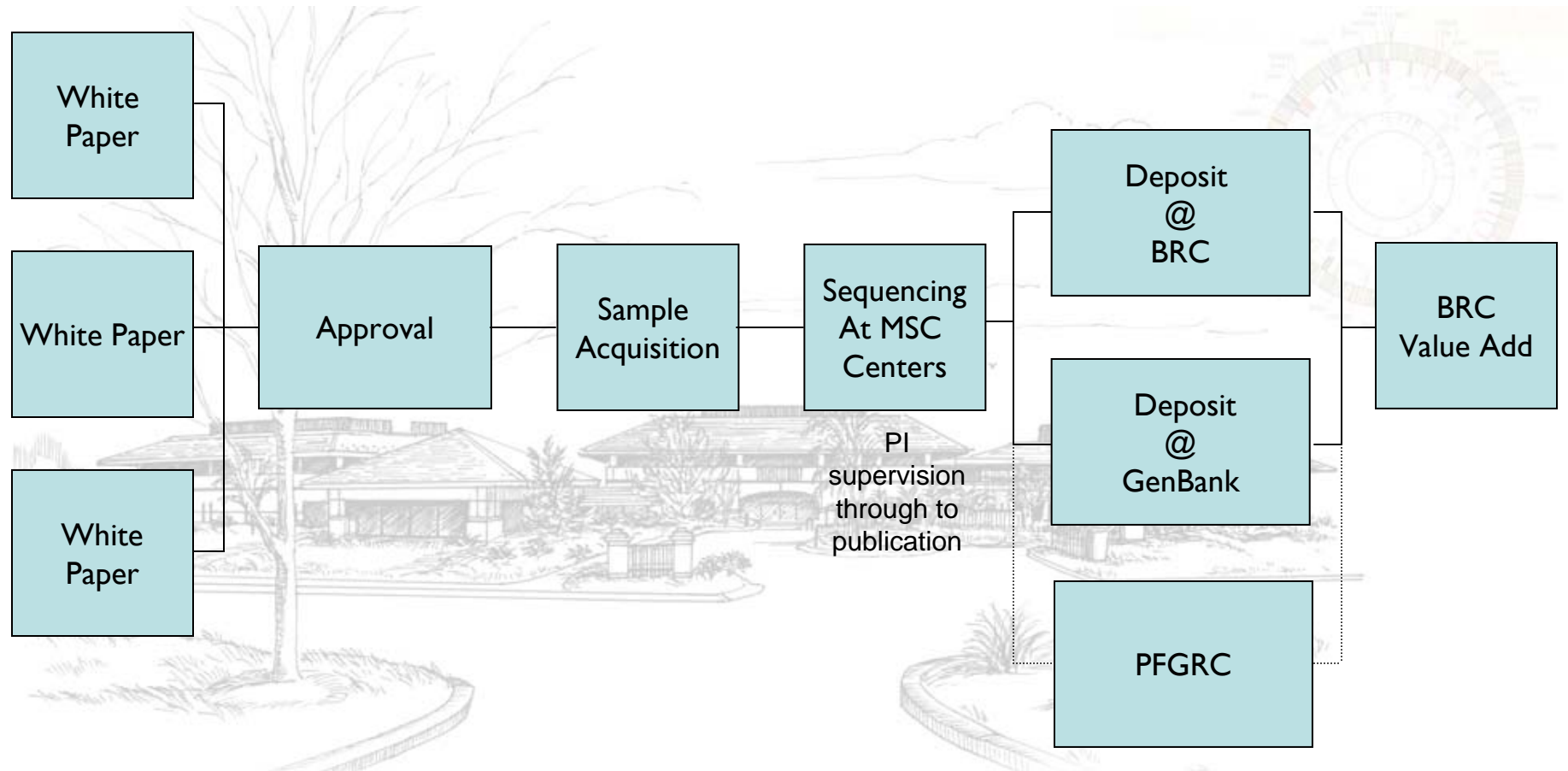
**Cloned
Genomic
Sequence**



Whole Genomic Sequence




Genome Work and Data Flow Across the TIGR/NIAID Resource Centers



Data Release & Reporting

<http://msc.tigr.org/status.shtml>



TIGR > MSC > Production Status

Production Status

TIGR's current MSC effort includes the following genome projects. Species/strain names link to the appropriate MSC site, if any, and other links lead to GenBank records. A 'p' indicates submission is in progress, and an 's' indicates the data has been submitted but is not yet available.

Recent status changes are listed on the [What's New](#) page. For a list of genome projects with individual MSC web sites, please see the [Genome Projects](#) page.

The projects on this page are organized into the following categories:

- [Disease Vectors](#)
- [NIAID Category A Pathogens](#)
- [NIAID Category B Pathogens](#)
- [NIAID Category C Pathogens](#)
- [Other Pathogens](#)
- [Related Species](#)

Taxon	Trace Archive	WGS	Assembly Archive	#Traces
Disease Vectors				
Aedes aegypti	X	I	W	7,747,816
Culex pipiens				
quinquefasciatus	X	I		3,596,685
Ixodes scapularis	X	I		6,293,858
NIAID Category A Pathogens				
Bacillus anthracis				
Tsiankovskii-I				
Yersinia pestis				
Angola	X	I	W	73,928
B42003004	X	I		72,991
E1979001	X	I		66,614
F1991016	X	I		65,198
IP275	X	I	W A	59,811

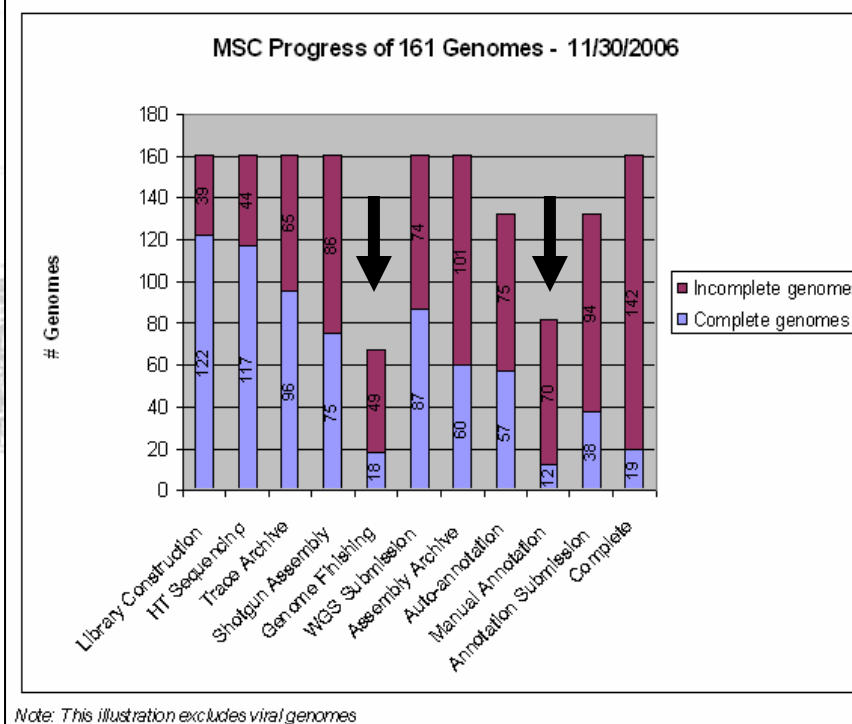
Submissions

- **Chromatograms:** daily or regular intervals based on project
- **Assembly files:** $\geq 3X$ assembly, 45 days after completion
- **Auto-annotation:** $\geq 5X$ immediately after completion
- **Manual Annotation:** After sequencing/finishing is completed

Current Status & Completion

#	ID	Name	Strains	Coverage / Closure	Auto Annotation	Manual Annotation
Insect Vectors of Disease						
1	8010B	<i>Aedes</i>	1	n(8X)	y	n
2	8010T	<i>Culex pipiens</i>	1	n(4X)	y	n
3	8010O	<i>Ixodes scapularis</i>	1	n(3X)	y	n
4	8010N	<i>R. communis</i>	1	n(4X)	n	n
Parasites						
5	8010Z	<i>Cryptosporidium muris</i>	1	n(8X)	y	y
6	8010M	<i>Entamoeba</i>	3	y(1)	y	y
7	8010F	<i>Plasmodium vivax</i>	1	y 9Mb	y	y
8	8010G	<i>Trichomonas</i> Genome	1	n(6X)	y	n
9	8010P	<i>Toxoplasma. Gondii</i>	2	n	y	y
Virus						
10	8010E	RNA Lab and Influenza	5400	y	n	n
11	8010K	Coronavirus Expansion	105	y	n	n
Fungi						
12	8010J	Comparative Aspergillus project	2	n	y	y
13	8010V	<i>Penicillium</i> spp	2	n	n	n
Bacteria						
14	8010R	<i>B. cereus</i>	10	y(5) 27.5 Mb	y	y
15	8010S	<i>Bartonella bacilliformis</i>	1	y y(17)	y	y
16	8010aa	<i>Borrelia</i>	17	0.69Mb/genome	y	y
17	8010A	<i>Burkholderia mallei</i>	11	y(2)	n	n
18	8010D	<i>Burkholderia pseudomallei</i>	9	y(3) 14.6Mb	n	n
19	8010Q	<i>Campylobacter</i>	10	y(5) 9Mb	y	n
20	8010L	<i>Coxiella</i>	6	y(6) 10Mb	y	y(2)
21	8010H	Pathogenic <i>E. coli</i> and <i>Shigella</i>	11	y(4)	y	n
22	8010ab	<i>E.coli</i> SECEC SMS3-5	1	y	y	y
23	8010U	<i>Pseudomonas aeruginosa PA7</i>	1	y	n	n
24	8010Y	<i>Salmonella</i> Strains	17	y(9) 44.1Mb	y	y(3)
25	8010X	<i>Ureaplasma</i>	20	y(2) 2Mb	y	y(2)
26	8010C	<i>Vibrio</i> strains	17	y(1)	y	y(1)
27	8010I	<i>Yersinia pestis</i>	9	y(3) 9.2 Mb	y	n
28	8010W	454 Data Assembly R&D	n/a	n/a	n/a	n/a

Contract Milestone Deliverables in the Aggregate	Due Date	% Completed	Deliverable Status
• Sequence submissions	Sept 2008	60%	On schedule
• WGS submissions	Sept 2008	54%	On schedule
• Closure and WGS updates	Sept 2008	27%	Behind schedule
• Annotation submissions	Sept 2008	29%	Behind schedule



New MSC Projects

Organism	Strains	Relevant BRC
Methicillin-resistant <i>Staphylococcus</i> plasmids	100 plasmids	NMPDR
<i>Vibrio parahaemolyticus</i>	7 strains	NMPDR
<i>Vibrio cholerae</i> – ICE	1 strain	NMPDR
<i>Escherichia coli</i> O157:H7	8 strains	ERIC
<i>Neisseria meningitidis</i>	Multiple genes	Not applicable
Murine coronavirus	10 strains	PATRIC

Selection Criteria

- Medical importance
- Status of other projects on the same organism
- Lab strain vs. clinical isolate
- Availability of experimental tools
- Community of users
- How well the strain represents the species

Sample Acquisition

- DNA / RNA
- Cell pellet

Genomes Common to NIAID/TIGR Resource Centers

#	MSC Genomes	Strains	BRC Presence
Insect Vectors of Disease			
1	<u>Aedes</u>	1	VectorBase
2	<u>Culex pipiens</u>	1	VectorBase
3	<u>Ixodes scapularis</u>	1	VectorBase
4	<u>R. communis</u>	1	
Parasites			
5	<u>Cryptosporidium muris</u>	1	ApiDB
6	<u>Entamoeba</u>	3	Pathema (E.histolytica)
7	<u>Plasmodium vivax</u>	1	ApiDB
8	<u>Trichomonas Genome</u>	1	
9	<u>Toxoplasma. Gondii</u>	2	
Virus			
10	RNA Lab and Influenza	5400	BioHealthBase
11	Coronavirus Expansion	105	Patric
Fungi			
12	Comparative Aspergillus project	2	
13	<u>Penicillium spp</u>	2	
Bacteria			
14	<u>B. cereus</u>	10	
15	<u>Bartonella bacilliformis</u>	1	
16	<u>Borrelia</u>	17	
17	<u>Burkholderia mallei</u>	11	Pathema
	<u>Burkholderia pseudomallei</u>	9	Pathema
18	<u>Campylobacter</u>	10	NMPDR (C.jejuni)
19	<u>Coxiella</u>	6	Patric (C.burnetti)
20	<u>Pathogenic E. coli and Shigella</u>	11	Eric
21	<u>E.coli SECEC SMS3-5</u>	1	
22	<u>Pseudomonas aeruginosa</u>	1	
23	<u>PA7</u>	1	
24	<u>Salmonella Strains</u>	17	
25	<u>Ureaplasma</u>	20	
26	<u>Vibrio strains</u>	17	
27	<u>Yersinia pestis</u>	9	Eric

BRC Center	BRC Genome	MSC presence	PFGRC Arrays
ApiDB	<u>Cryptosporidium spp.</u>	✓	
	<u>Plasmodium spp.</u>	✓	✓
	<u>Toxoplasma gondii</u>	✓	
BioHealthBase	<u>Francisella tularensis</u>		✓
	<u>Giardia lamblia</u>		✓
	<u>Microsporidia</u>		
	<u>Ricinus communis</u>	✓	
	<u>Mycobacterium tuberculosis</u>		✓
ERIC	<u>Influenza Virus</u>	✓	
	<u>Yersinia pestis</u>	✓	✓
	<u>Diarrheagenic E. coli</u>	✓	
NMPDR	<u>Shigella</u>	✓	
	<u>Salmonella</u>	✓	✓
	<u>Staphylococcus aureus</u>		
Pathema	<u>Pathogenic Vibrios</u>	✓	✓
	<u>Listeria monocytogenes</u>		✓
	<u>Campylobacter jejuni</u>	✓	✓
	<u>Streptococcus pyogenes</u>		
	<u>Streptococcus pneumoniae</u>		✓
Patric	<u>Bacillus anthracis</u>		✓
	<u>Burkholderia mallei</u>	✓	✓
	<u>Burkholderia pseudomallei</u>	✓	✓
	<u>Clostridium perfringens</u>	✓	
	<u>Entamoeba histolytica</u>	✓	
VBRC	<u>Rickettsiae</u>		✓
	<u>Brucella</u>		
	<u>Coxiella burnetii</u>	✓	
	<u>Calicivirus</u>		
	<u>Hepatitis A Virus</u>		
VectorBase	<u>Rabies Virus</u>		
	<u>Coronavirus</u>	✓	✓
	<u>Variola major Virus</u>		
	<u>Arenavirus</u>		
	<u>Hanta Virus</u>		
	<u>Rift Valley Fever Virus</u>		
	<u>Ebola Virus</u>		
	<u>Marburg Virus</u>		
	<u>Dengue Virus</u>		
	<u>California encephalitis group Virus</u>		
	<u>Kyasanar forest disease Virus</u>		
	<u>Omsk hemorrhagic fever Virus</u>		
	<u>West Nile Virus</u>		
	<u>Alphavirus</u>		
	<u>Hantaan Virus</u>		
	<u>Puumala Virus</u>		
	<u>Crimean-Congo hemorrhagic fever Virus</u>		
	<u>Yellow fever Virus</u>		
	<u>Tick-borne Encephalitis</u>		
	<u>Nipah Virus</u>		
	<u>Equine morbillivirus</u>		
	<u>Anopheles gambiae</u>		
	<u>Aedes aegypti</u>	✓	
	<u>Anopheles gambiae</u>		
	<u>Culex pipiens</u>	✓	
	<u>Ixodes scapularis</u>	✓	

MSC Program



BRC Genomes



PFGRC Arrays



Microbial Sequencing Center



[TIGR](#) > [MSC](#) > [Genome Projects](#) > [Yersinia](#)

About Us
What's New
Production Status
Genome Projects
Yersinia
<ul style="list-style-type: none">GoalsIsolate SelectionStrainsGenBank DataReferences
Yersinia pestis
<ul style="list-style-type: none">AngolaB42003004E1979001F1991016IP275

Y. pestis str:
rhamnose a
Yersinia pse
Y. pestis str:
pseudotube
pestis strain
typical *Y. pe*
Yersinia, the
pestis. Ango
deletion affe



THE INSTITUTE FOR GENOMIC RESEARCH

MSC

Microbial Sequencing Center

| [Contact Us](#) | [TIGR](#) | [NIAID](#) |

[TIGR](#) > [MSC](#) > [Genome Projects](#) > [Aedes aegypti](#) > Released Data

About Us
What's New
Production Status
Genome Projects
<i>A. aegypti</i>
<ul style="list-style-type: none">Released DataGenBankSubmissionsBAC EndsTIGR Gene IndexSequencing InformationGenBank DataProgress by LibraryTimeline [PDF]Project Plan [PDF]
Publications
Resources
Data Release Policies
Team Members
>> NIAID MSC

Released Data

Updated 6/14/06

Aedes aegypti annotation Release 1.0

The NIAID Microbial Sequencing Centers are pleased to announce annotation Release 1.0 of the *Aedes aegypti* genome sequence. This annotation was produced jointly by [The Institute for Genomic Research](#) and [VectorBase](#) with support from [The Broad Institute of Harvard/MIT](#).

Release 1.0 has been deposited at GenBank under the accession version [AGE00000000](#). Long-term curation of the genome sequence and subsequent annotation updates will be the responsibility of [VectorBase](#).

Release 1.0 contains 15,419 high confidence protein-coding genes; alternatively spliced transcripts derived from 992 genes add an additional 1,370 proteins yielding a total of 16,789 predicted proteins. Our structural annotation (gene finding) strategy involved masking repetitive DNA sequences, followed by a combination of EST and protein alignments, trained gene prediction algorithms, and evidence-based gene predictions. Annotation release 1.0 was derived by comparing and merging gene sets generated independently by VectorBase and TIGR.

An autonaming pipeline used sequence similarity to PANTHER, Drosophila (Release 4.3) and UniProt databases as a basis for functional name assignment. Gene Ontology terms were computationally assigned to the genes based on Drosophila GO annotations via sequence homology, or Pfam domain matches above the trusted cutoff. A "supplement" of lower confidence gene predictions is also available at [VectorBase](#).

Aedes MSC and VectorBase Project Plan

The National Institute of Allergy and Infectious Diseases, National Institutes of Health has funded the *Aedes aegypti* genome project through its Microbial Sequencing Centers (MSCs) at [The Institute for Genomic Research \(TIGR\)](#) and the [Broad Institute](#). The MSCs are responsible for genome sequencing, assembly, and annotation of gene structure and function, with the goal of rapid release of each of these data sets to the scientific community. Once released, the complete sequence and annotation of the *Aedes* genome will permanently reside at a third [NIAD-sponsored entity](#), [VectorBase](#), which is a Bioinformatics Resource Center (BRC) at the University Of Notre Dame. Delivery of these data will also require coordination with NCBI-GenBank. Given the mutual interests of these organizations, the most effective approach to the initial release of *Aedes* genomic data will be to work in close collaboration to produce an initial set of annotation, refine and improve the pipelines resident at each of the centers, and to generate data to use in the analysis and publication of *Aedes*.

We are confident that the combined annotation efforts of the MSCs and VectorBase produce unified, high-quality *Aedes* annotation for release to the scientific community. This document is a detail project plan intended to describe the complementary activities planned among the three NIAD-funded centers.

See the accompanying document [Timeline.pdf](#) for the task timeline overview.

A summary of the activities for the MSC *Aedes* project listed by institution

Institution: Broad			
Task	Duration	Start	Finish
Stage 0: Establish Communication Mechanisms Among TIGR, Broad and VectorBase	5 weeks	9/12/05	10/14/05
Develop Assembly Release Strategy	3 weeks	9/12/05	
Conference Call	1 day	9/27/05	
Review and Acceptance of Project Plan	2.4 weeks	9/22/05	
Stage 1: Annotation Preparation	4 weeks	9/19/05	
Release August Assembly	5 days	10/03/05	
Release useful datasets	5 days	10/10/05	
Run and Evaluate Computes	3 weeks	9/26/05	
Stage 2: Production Annotation Gene Structure	8 weeks	10/17/05	

Institution: VectorBase			
Task	Duration	Start	Finish
Stage 0: Establish Communication Mechanisms Among TIGR, Broad and VectorBase	5 weeks	9/12/05	10/14/05
Conference Call	1 day	9/27/05	
Review and Acceptance of Project Plan	2.4 weeks	9/22/05	
Stage 1: Annotation Preparation	4 weeks	9/19/05	10/14/05
Release useful datasets	5 days	10/10/05	10/14/05
Stage 2: Production Annotation Gene Structure	8 weeks	10/17/05	12/09/05
VectorBase: gene structure annotation	8 weeks	10/17/05	12/09/05

A summary of the Stages of completion for the MSC *Aedes* project

Stage 0: Establish Communication Mechanisms among TIGR, Broad and VectorBase (5 weeks)

Stage 1: Annotation Preparation (4 weeks)

Stage 2: Production Annotation Gene Structure (8 weeks)

Stage 3: Evaluation of Data (4 weeks)

Stage 4: Data Generation Gene Structure v1.0 (4 weeks)

Stage 5: Production Annotation Functional Computes (3 weeks)

Stage 6: Genbank Submission v1.0 (2 weeks)

Stage 7: Genbank Processing to Release (4 weeks)

Stage 8: Manuscript Preparation (29 weeks)

ID	Task Name	Duration	Start	Finish	Predecess	Notes	Resource Names	Successors
1	Stage 0: Establish Communication Mechanisms Among TIGR, Broad and VectorBase	25 days	Mon 9/12/05	Fri 10/14/05		ication, such as conference calls or email ...	Broad,TIGR,VectorBase	12
2	Develop Assembly Release Strategy	15 days	Mon 9/12/05	Fri 9/30/05		3 the primary annotation to the Broad assembly. ...	Broad,TIGR	9
3	Develop Draft Project Plan for Annotation and Analysis	8 days	Mon 9/12/05	Wed 9/21/05		circulation. The timeframe represented in the d...	TIGR	6,5,11
4	Define Annotation Data types and Metrics for Evaluating Gene Sets	15 days	Mon 9/12/05	Fri 9/30/05		Annotation Data types and File Formats ...	TIGR	
5	Conference Call	0 days	Tue 9/27/05	Tue 9/27/05	3	process edits and arrive at a ratified document.	Broad,TIGR,VectorBase	
6	Review and Acceptance of Project Plan	12 days	Thu 9/22/05	Fri 10/7/05	3	eviewed by NIAD and a final version distributed.	Broad,TIGR,VectorBase,NIAD	29
7	Stage 1: Annotation Preparation	20 days	Mon 9/19/05	Fri 10/14/05		uction Annotation will occur at this stage. ...	Broad,TIGR,VectorBase	12
8	Create Central Repository (CR)	10 days	Mon 9/19/05	Fri 9/30/05		anged between the Broad, TIGR and VectorBase.	TIGR	10
9	Release August Assembly	5 days	Mon 10/03/05	Fri 10/07/05	2	n the Broad on a project submission date.	Broad	10
10	Release useful datasets	5 days	Mon 10/10/05	Fri 10/14/05	9,8	leased to the scientific community. The informat...	TIGR,Broad,VectorBase	
11	Define the 10mb region for manual annotation	8 days	Thu 9/22/05	Mon 10/03/05	3	reloped by TIGR and VectorBase will be identified.	Broad	
12	Stage 2: Production Annotation Gene Structure	40 days?	Mon 10/17/05	Fri 12/9/05	1,7	How us to perform comprehensive evalua...	TIGR,Broad,VectorBase	29,47
13	TIGR gene structure annotation	40 days	Mon 10/17/05	Fri 12/9/05		pertain to TIGR's gene structure annotation.	TIGR	
14	Run autpeline	4 wks	Mon 10/17/05	Fri 11/11/05		g), to the larger and more complex genomes of pl...	TIGR	15
15	Submit TIGR 0.5 in Central Repository	0 days	Fri 11/11/05	Fri 11/11/05	14	will be deposited in the Central Repository FTP site.	TIGR	16
16	Iterative improvement of gene set	15 days	Mon 11/14/05	Fri 12/2/05	15	the individual gene prediction programs. ...	TIGR	
17	Verify against Anopheles and Drosophila	5 days	Mon 11/14/05	Fri 11/18/05		compared to our predicted gene set to identify in...	TIGR	18
18	Locate genes in introns of others	5 days	Mon 11/21/05	Fri 11/25/05	17	nes were identified. Therefore, we will extract ...	TIGR	19
19	EST Data Incorporation	5 days	Mon 11/28/05	Fri 12/2/05	18	EST alignments and the gene structures, and c...	TIGR	20
20	Quality Control	1 wk	Mon 12/05/05	Fri 12/09/05	19	editing of existing gene models, others will resul...	TIGR	21
21	Submit 0.5.1 in Central Repository	0 days	Fri 12/09/05	Fri 12/09/05	20	l will be deposited in Central Repository FTP site.	TIGR	
22	VectorBase: gene structure annotation	40 days	Mon 10/17/05	Fri 12/9/05		erbrates and also Anopheles gambiae. ...	VectorBase	
23	autpeline	4 wks	Mon 10/17/05	Fri 11/11/05		rganisms. There are three main sources genes: ...	VectorBase	24
24	make VB 0.5 available via web	4 wks	Mon 11/14/05	Fri 12/09/05	23	lable to community for BLAST, SSAHA access...	VectorBase	25
25	Submit 0.5 or 0.5+ into Central Repository	0 days	Fri 12/09/05	Fri 12/09/05	24	will be deposited in Central Repository FTP site...	VectorBase	
26	Broad: 10 Mb annotation	40 days?	Mon 10/17/05	Fri 12/9/05		b). This will provide a complementary dat...	Broad	
27	Annotate 10 Mb region	8 wks?	Mon 10/17/05	Fri 12/09/05		Station for the 10Mb region is requested here. 	Broad	28
28	Submit 10 Mb annotation in Central Repository	0 days	Fri 12/09/05	Fri 12/09/05	27	l, will be deposited in Central Repository FTP site.	Broad	
29	Stage 3: Evaluation of Data	20 days?	Mon 12/12/05	Fri 1/6/06	12,6	will be deposited into Genbank in Stage 5...	Broad,TIGR,VectorBase	40,33
30	Implement work in Proposal for evaluating and comparing gene sets	15 days	Mon 12/12/05	Fri 12/30/05		The tasks required for data evaluation include...	TIGR	31,32
31	define merge strategy based on evaluation	5 days	Mon 12/26	Fri 1/6/06	30	decision on the composition of the final gene set.	TIGR,Broad,VectorBase	

Would better BRC-MS C coordination matter?

in support of diagnostics, therapeutics, vaccines, mechanisms, virulence factors etc.

“To call in a statistician after the experiment is done may be no more than asking him to perform a postmortem examination: he may be able to say what the experiment died of.”
(R.A. Fisher)

- “Faith-based” bottom-up initiatives like the resource centers actually work
- Requirements inputs from expert community and applications developers to MSCs and BRCs to focus efforts and outputs with well-posed problems/questions
- Project management plans to push/pull MSC genomes and MSC generated data and analyses to their corresponding BRCs
- Are MSCs and BRCs separated by a common language and uncommon metrics?
- Collective future of Resource Centers like MSCs and BRCs depends on demonstrating value

Thanks

- NIAID
- TIGR Service Team members of
 - The Joint Technology Center (JTC) – Sequencing
 - Biotechnology Core Services (BCLS & BCIS) – Genome finishing
 - Informatics (IFX) - Assembly, Annotation & Data Submissions
- **Ishwar Chandramouliswaran, MSC Project Manager**